

Electronic Letters on Computer Vision and Image Analysis 5(3):96-104, 2005

Hand Tracking and Gesture Recognition for Human-Computer Interaction

Cristina Manresa, Javier Varona, Ramon Mas and Francisco J. Perales

*Unidad de Gráficos y Visión por Computador
Departamento de Matemáticas e Informática
Universitat de les Illes Balears
Edificio Anselm Turmeda, Ctra. Valldemossa km 7.5
07122 – Palma de Mallorca - España*

Received 4 February 2005; accepted 18 May 2005

Abstract

The proposed work is part of a project that aims at the control of a videogame based on hand gesture recognition. This goal implies the restriction of real-time response and the use of unconstrained environments. In this paper we present a new algorithm to track and recognise hand gestures for interacting with a videogame. This algorithm is based on three main steps: hand segmentation, hand tracking and gesture recognition from hand features. For the hand segmentation step we use the colour cue due to the characteristic colour values of human skin, its invariant properties and its computational simplicity. To prevent errors from hand segmentation we add the hand tracking as a second step. Tracking is performed assuming a constant velocity model and using a pixel labeling approach. From the tracking process we extract several hand features that are fed into a finite state classifier which identifies the hand configuration. The hand can be classified into one of the four gesture classes or one of the four different movement directions. Finally, the system's performance is evaluated by showing the usability of the algorithm in a videogame environment.

Key Words: Hand Tracking, Gesture Recognition, Human-Computer Interaction, Perceptual User Interfaces.

1 Introduction

Nowadays, the majority of human-computer interaction (HCI) is based on mechanical devices such as keyboards, mice, joysticks or gamepads. In recent years there has been a growing interest in methods based on computational vision due to its ability to recognise human gestures in a natural way [1]. These methods use the images acquired from a camera or from a stereo pair of cameras as input. The main goal of these algorithms is to measure the hand configuration at each time instant.

To facilitate this process many gesture recognition applications resort to the use of uniquely coloured gloves or markers on hands or fingers [2]. In addition, using a controlled background makes it possible to locate the hand efficiently, even in real-time [3]. These two conditions impose restrictions on the user and on the interface setup. We have specifically avoided solutions that require coloured gloves or markers and a

Correspondence to: <cristina.manresa@uib.es>

Recommended for acceptance by <Perales F., Draper B.>

ELCVIA ISSN: 1577-5097

Published by Computer Vision Center / Universitat Autònoma de Barcelona, Barcelona, Spain

controlled background because of the initial requirements of our application. It must work for different people, without any complement on them and also for unpredictable backgrounds.

Our application uses images from a low-cost web camera placed in front of the work area, where the recognised gestures act as the input for a computer 3D videogame. The players, rather than pressing buttons, must use different hand gestures that our application should recognise. This fact, increases the complexity since the response time must be very fast. Users should not appreciate a significant delay between the instant they perform a gesture or motion and the instant the computer responds. Therefore, the algorithm must provide real-time performance for a conventional processor. Most of the known hand tracking and recognition algorithms do not meet this requirement and are inappropriate for visual interface. For instance, particle filtering-based algorithms can maintain multiple hypotheses at the same time to robustly track the hands but they need high computational demands [4]. Recently, several contributions for reducing the complexity of particle filters have been presented, for example, using a deterministic process to help the random search [5]. Also in [6], we can see a multi-scale colour feature for representing hand shape and particle filtering that combines shape and colour cues in a hierarchical model. The system has been fully tested and seems robust and stable. To our knowledge the system runs at about 10frames/second and does not consider several hand states. However, these algorithms only work in real-time for a reduced size hand and in our application, the hand fills most of the image. In [7], shape reconstruction is quite precise, a high DOF model is considered, and in order to avoid self-occlusions infrared orthogonal cameras are used. The authors propose to apply this technique using a colour skin segmentation algorithm.

In this paper we propose a real-time non-invasive hand tracking and gesture recognition system. In the next sections we explain our method which is divided in three main steps. The first step is hand segmentation, the image region that contains the hand has to be located. In this process, the use of the shape cue is possible, but they vary greatly during the natural hand motion[8]. Therefore, we choose skin-colour as the hand feature. The skin-colour is a distinctive cue of hands and it is invariant to scale and rotation. The next step is to track the position and orientation of the hand to prevent errors in the segmentation phase. We use a pixel-based tracking for the temporal update of the hand state. In the last step we use the estimated hand state to extract several hand features to define a deterministic process of gesture recognition. Finally, we present the system's performance evaluation results that prove that our method works well in unconstrained environments and for several users.

2 Hand Segmentation Criteria

The hand must be located in the image and segmented from the background before recognition. Colour is the selected cue because of its computational simplicity, its invariant properties regarding to the hand shape configurations and due to the human skin-colour characteristic values. Also, the assumption that colour can be used as a cue to detect faces and hands has been proved useful in several publications [9,10]. For our application, the hand segmentation has been carried out using a low computational cost method that performs well in real time. The method is based on a probabilistic model of the skin-colour pixels distribution. Then, it is necessary to model the skin-colour of the user's hand. The user places part of his hand in a learning square as shown in Fig. 1. The pixels restricted in this area will be used for model learning. Next, the selected pixels are transformed from the RGB-space to the HSL-space and the chroma information is taken: hue and saturation.



Figure 1: Application interface and skin-colour learning square.

We have encountered two problems in this step that have been solved in a pre-processing phase. The first one is that human skin hue values are very near to red colour, that is, their value is very close to 2π radians, so it is difficult to learn the distribution due to the hue angular nature that can produce samples on both limits. To solve this inconvenience the hue values are rotated π radians. The second problem in using HSL-space appears when the saturation values are close to 0, because then the hue is unstable and can cause false detections. This can be avoided discarding saturation values near 0.

Once the pre-processing phase has finished, the hue, H , and saturation, S , values for each selected pixel are used to infer the model, that is, $\bar{x} = (\bar{x}_1, \dots, \bar{x}_n)$, where n is the number of samples and a sample is $\bar{x}_i = (H_i, S_i)$. A Gaussian model is chosen to represent the skin-colour probability density function. The values for the parameters of the Gaussian model (mean, \bar{x} , and covariance matrix, Σ) are computed from the sample set using standard maximum likelihood methods [11]. Once they are found, the probability that a new pixel, $\bar{x} = (H, S)$, is skin can be calculated as

$$P(\bar{x} \text{ is skin}) = \frac{1}{\sqrt{(2\pi)^2 |\Sigma|}} e^{(-\frac{1}{2}(\bar{x} - \bar{x}) \Sigma^{-1} (\bar{x} - \bar{x})^T)} \quad (1)$$

Finally, we obtain the blob representation of the hand by applying a connected components algorithm to the probability image, which groups pixels into the same blob. The system is robust to background changes and low light conditions. If the system gets lost, you can initialise it again by going to the hand start state. Fig. 2 shows the blob contours found by the algorithm for different environment conditions where the system has been tested.

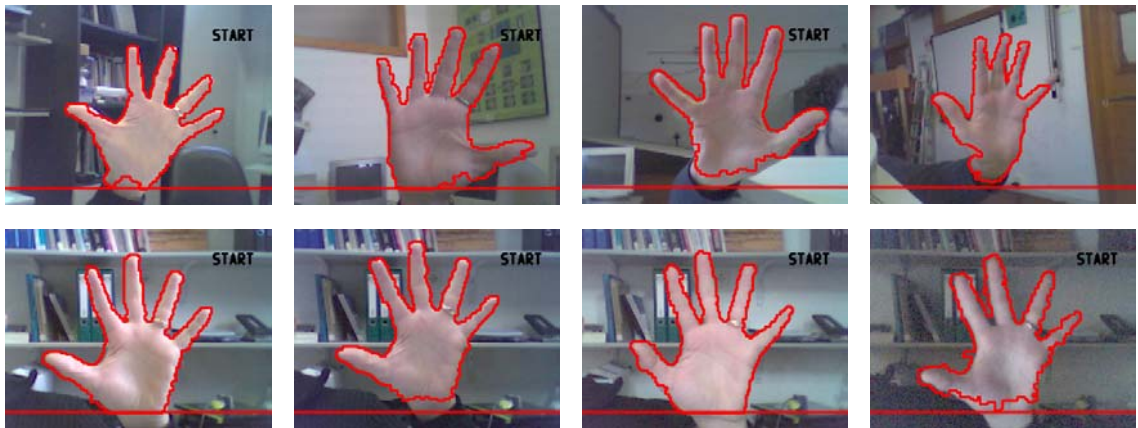


Figure 2: Hand contours for different backgrounds (1st row) and different light conditions (2nd row).

3 Tracking Procedure

USB cameras are known for the low quality images they produce. This fact can cause errors in the hand segmentation process. In order to make the application robust to these segmentation errors we add a tracking algorithm. This algorithm tries to maintain and propagate the hand state over time.

We represent the hand state in time t , $\vec{s}(t)$, by means of a vector, $\vec{s}(t) = (\vec{p}(t), \vec{w}(t), \alpha(t))$, where $\vec{p} = (p_x, p_y)$ is the hand position in the 2D image, the hand size is represented by $\vec{w} = (w, h)$, where w is the hand width and h is the hand height in pixels, and, finally, α is the hand's angle in the 2D image plane. First, from the hand state in time t we built a hypothesis of the hand state, $\vec{h} = (\vec{p}(t+1), \vec{w}(t), \alpha(t))$, for time $t+1$ applying a simple second-order autoregressive process to the position component

$$\vec{p}(t+1) - \vec{p}(t) = \vec{p}(t) - \vec{p}(t-1). \quad (2)$$

Equation (2) expresses a dynamical model of constant velocity. Next, if we assume that at time t , M blobs have been detected, $B = \{b_1, \dots, b_j, \dots, b_M\}$, where each blob b_j corresponds to a set of connected skin-colour pixels, the tracking process has to set the relation between the hand hypothesis, \vec{h} , and the observations, b_j , over time.

In order to cope with this problem, we define an approximation to the distance from the image pixel, $\vec{x} = (x, y)$, to the hypothesis \vec{h} . First, we normalize the image pixel coordinates

$$\vec{n} = R \cdot (\vec{x} - \vec{p}(t+1)), \quad (3)$$

where R is a standard 2D rotation matrix about the origin, α is the rotation angle, and $\vec{n} = (n_x, n_y)$ are the normalized pixel coordinates. Then, we can find the crossing point, $\vec{c} = (c_x, c_y)$, between the hand hypothesis ellipse and the normalized image pixel as follows

$$\begin{aligned} c_x &= w \cdot \cos \mathcal{G} \\ c_y &= h \cdot \sin \mathcal{G} \end{aligned} \quad (4)$$

where \mathcal{G} is the angle between the normalized image pixel and the hand hypothesis. Finally, the distance from an image pixel to the hand hypothesis is

$$d(\vec{x}, h) = \|\vec{n}\| - \|\vec{c}\|. \quad (5)$$

This distance can be seen as the approximation of the distance from a point in the 2D space to a normalized ellipse (normalized means centered in origin and not rotated). From the distance definition of (5) it turns out that its value is equal or less than 0 if \vec{x} is inside the hypothesis h , and greater than 0 if it is outside. Therefore, considering the hand hypothesis h and a point \vec{x} belonging to a blob b , if the distance is equal or less than 0, we conclude that the blob b supports the existence of the hypothesis h and it is selected to represent the new hand state. This tracking process could also detect the presence or the absence of the hand in the image [12].

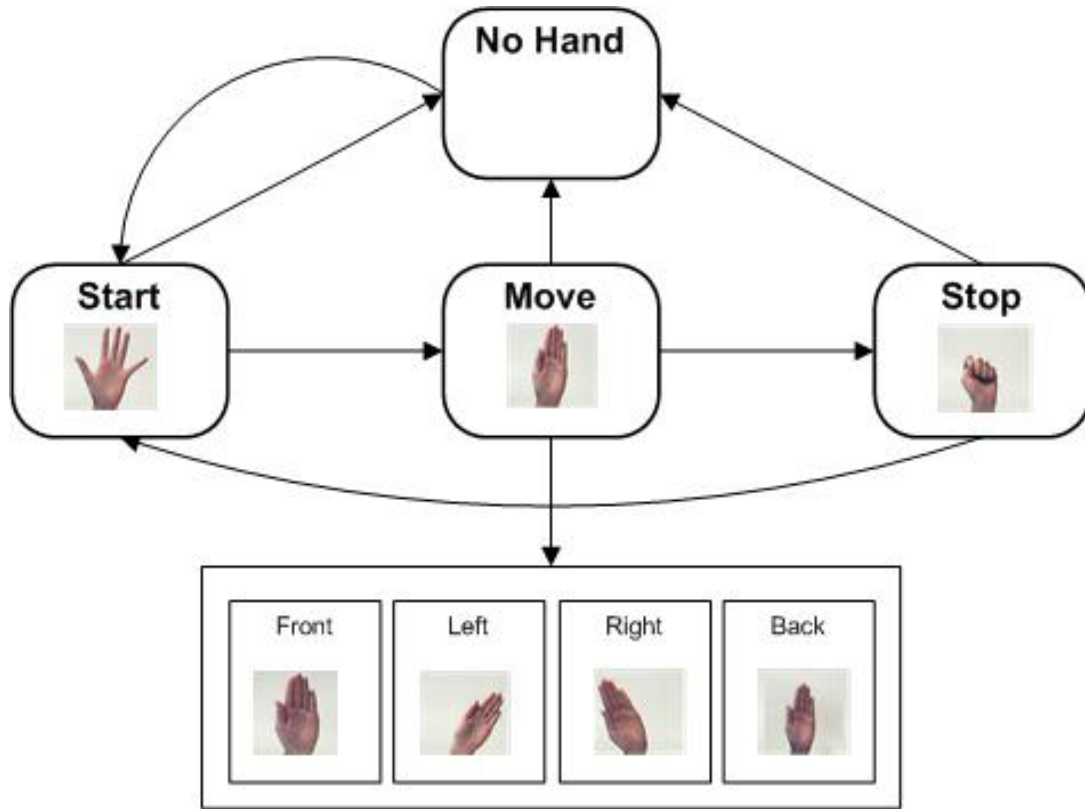


Figure 3: Gesture alphabet and valid gesture transitions.

4 Gesture Recognition

Our gesture alphabet consists in four hand gestures and four hand directions in order to fulfil the application's requirements. The hand gestures correspond to a fully opened hand (with separated fingers), an opened hand with fingers together, a fist and the last gesture appears when the hand is not visible, in part or completely, in the camera's field of view. These gestures are defined as *Start*, *Move*, *Stop* and the *No-Hand* gesture respectively. Also, when the user is in the *Move* gesture, he can carry out *Left*, *Right*, *Front* and *Back* movements. For the *Left* and *Right* movements, the user will rotate his wrist to the left or right. For the *Front* and *Back* movements, the hand will get closer to or further from the camera. Finally, the valid hand gesture transitions that the user can carry out are defined in Fig. 3.

The process of gesture recognition starts when the user's hand is placed in front of the camera's field of view and the hand is in the *Start* gesture, that is, the hand is fully opened with separated fingers. In order to avoid fast hand gesture changes that were not intended, every change should be kept fixed for 5 frames, if not the hand gesture does not change from the previous recognised gesture.

To achieve this gesture recognition, we use the hand state estimated in the tracking process, that is, $\vec{s} = (\vec{p}, \vec{w}, \alpha)$. This state can be viewed as an ellipse approximation of the hand where $\vec{p} = (p_x, p_y)$ is the ellipse center and $\vec{w} = (w, h)$ is the size of the ellipse in pixels. To facilitate the process we define the major axis length as M and the minor axis length as m . In addition, we compute the hand's blob contour and its corresponding convex hull using standard computer vision techniques. From the hand's contour and the hand's convex hull we can calculate a sequence of contour points between two consecutive convex hull vertices. This sequence forms the so-called convexity defect (i.e., a finger concavity) and it is possible to compute the depth of the i th-convexity defect, d_i . From these depths it is possible to compute the depth

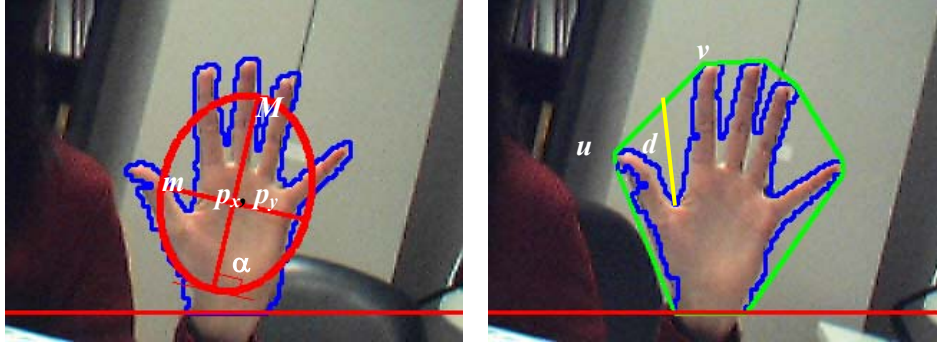


Figure 4: Extracted features for the hand gesture recognition. In the right image, u_i and v_i indicate the start and the end points of the i th-convexity defect, the depth, d_i , is the distance from the farthestmost point of the convexity defect to the convex hull segment.

average, \bar{d} , as a global hand feature, see (6), where n is the total number of convexity defects in the hand's contour, see Fig. 4.

$$\bar{d} = \frac{1}{n} \sum_{i=0..n} d_i. \quad (6)$$

The first step of the gesture recognition process is to model the *Start* gesture. The average of the depths of the convexity defects of an opened hand with separated fingers is larger than in an open hand with no separated fingers or in a fist. This feature is used for differentiating the next hand gesture transitions: from *Stop* to *Start*; from *Start* to *Move*; and from *No-Hand* to *Start*. However, first it is necessary to compute the *Start* gesture feature, T_{start} . Once the user is correctly placed in the camera's field of view with the hand widely opened the skin-colour learning process is initiated. The system also computes the *Start* gesture feature for the n first frames,

$$T_{\text{start}} = \frac{\frac{1}{n} \sum_{t=0..n} d(t)}{2}. \quad (7)$$

Once the *Start* gesture is identified, the most probable valid gesture change is the *Move* gesture. Therefore, if the current hand depth is less than T_{start} the system goes to the *Move* hand gesture. If the current hand gesture is *Move* the hand directions will be enabled: *Front*, *Back*, *Left* and *Right*.

If the user does not want to move in any direction, he should set his hand in the *Move* state. The first time that the *Move* gesture appears, the system computes the *Move* gesture feature, T_{move} , that is an average of the approximated area of the hand for n consecutive frames,

$$T_{\text{move}} = \frac{1}{n} \sum_{t=0..n} M(t) \cdot m(t). \quad (8)$$

In order to recognize the *Left* and *Right* directions, the calculated angle of the fitted ellipse is used. To prevent non desired jitter effects in orientation, we introduce a predefined constant T_{jitter} . Then, if the angle

of the ellipse that circumscribes the hand, α , satisfies $\alpha > T_{\text{jitter}}$, *Left* orientation will be set. If the angle of the ellipse that circumscribes the hand, α , satisfies $\alpha < -T_{\text{jitter}}$, *Right* orientation will be set.

In order to control the *Front* and *Back* orientations and to return to the *Move* gesture the hand must not be rotated and the *Move* gesture feature is used to differentiate these movements. If $T_{\text{move}} \cdot C_{\text{front}} < M \cdot m$ succeeds the hand orientation will be *Front*. The *Back* orientation will be achieved if $C_{\text{back}} > m/M$.

The *Stop* gesture will be recognised using the ellipse's axis. When the hand is in a fist, the fitted ellipse is almost like a circle and m and M are practically the same, that is, when $M - m < C_{\text{stop}} \cdot C_{\text{front}}, C_{\text{back}}$ and C_{stop} are predefined constants established during the algorithm performance evaluation. Finally, the *No-Hand* state will appear when the system does not detect the hand, the size of the detected hand is not large enough or when the hand is in the limits of the camera's field of view. The next possible hand state will be the *Start* gesture and it will be detected using the transition procedure from *Stop* to *Start* explained earlier on.

Some examples of gesture transitions and the recognised gesture results can be seen in Fig. 5. These examples are chosen to show the algorithm robustness for different lighting conditions, hand configurations and users. We realize that a correct learning of the skin-colour is very important. If not, some problems with the detection and the gesture recognition can be encountered. One of the main problems with the use of the application is the hand control, maintaining the hand in the camera's field of view and without touching the limits of the capture area. This problem has been shown to disappear with user's training.

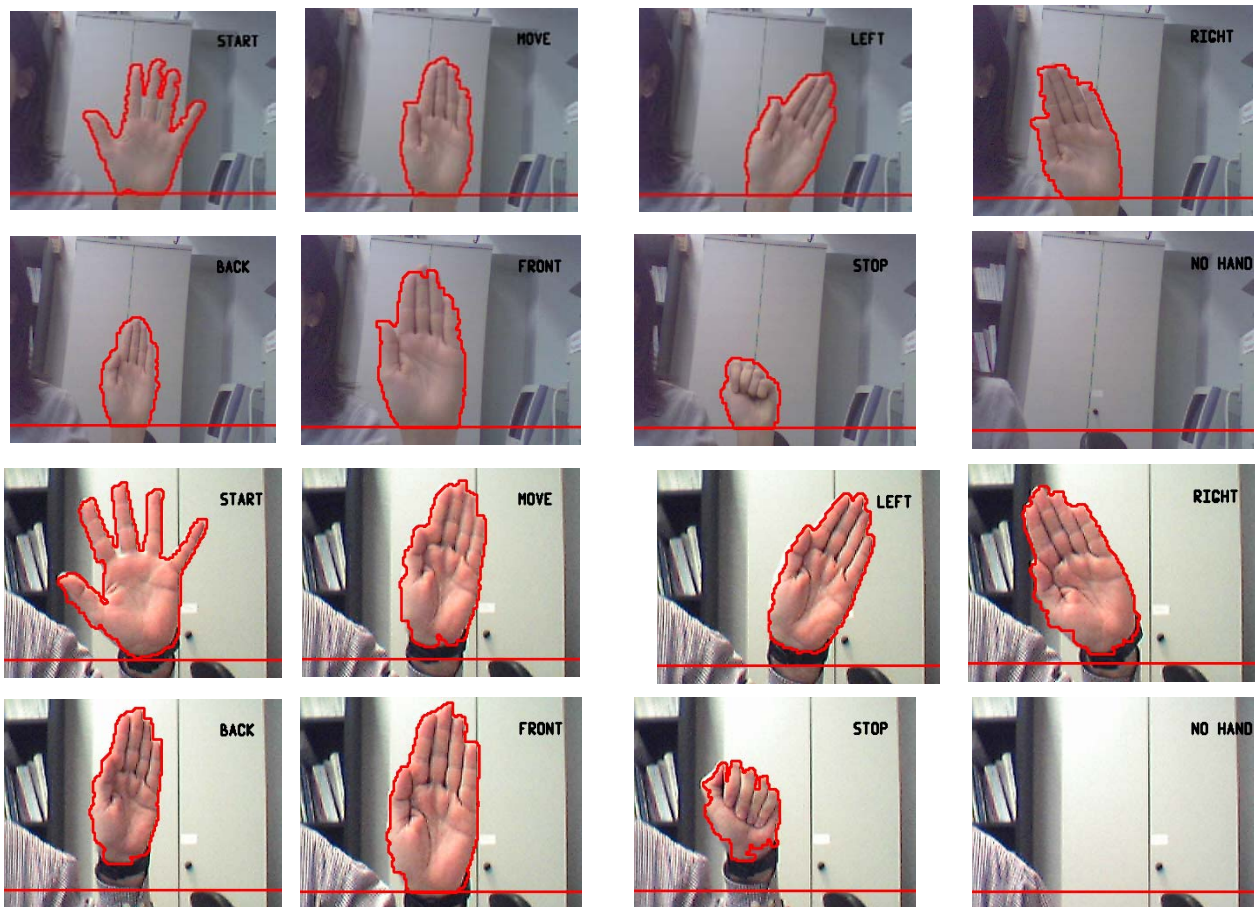


Figure 5: Gesture recognition examples for different lighting conditions, users and hand configurations.

5 System's performance evaluation

In this section we show the accuracy of our hand tracking and gesture recognition algorithm. The application has been implemented in Visual C++ using the OpenCV libraries [13]. The application has been tested on a Pentium IV running at 1.8 GHz. The images have been captured using a Logitech Messenger WebCam with USB connection. The camera provides 320x240 images at a capture and processing rate of 30 frames per second.

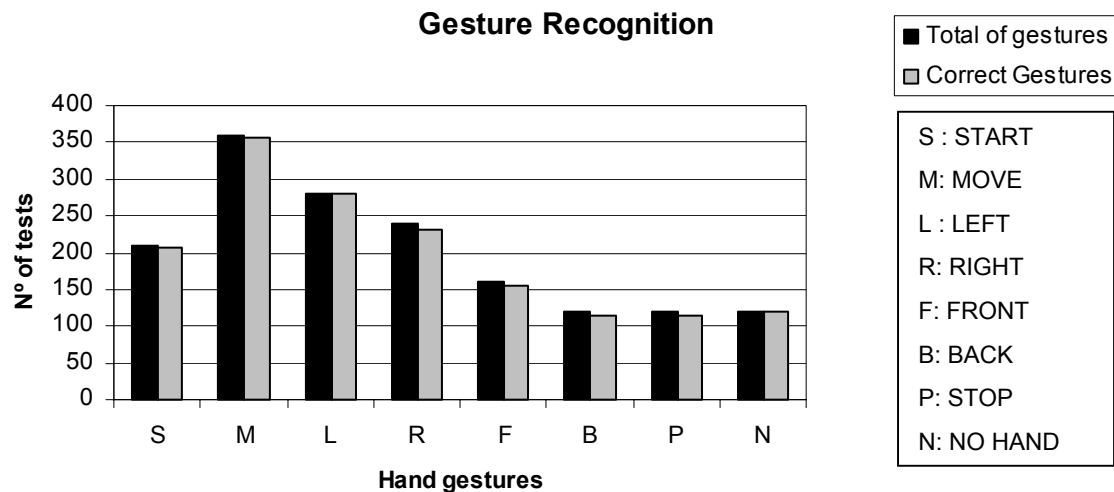


Figure 6: System's performance evaluation results.

For the performance evaluation of the hand tracking and gesture recognition, the system has been tested on a set of 40 users. Each user has performed a predefined set of 40 gestures and therefore we have 1600 gestures to evaluate the application results. It is natural to think that the system's accuracy will be measured controlling the performance of the desired user movements for managing the videogame. This sequence included all the application's possible states and transitions. Figure 6 shows the performance evaluation results. These results are represented using a graph with the application states, such as Start or Move, as columns and the number of appearances of the gesture as rows. The columns are paired for each gesture: the first column is the number of tests of the gesture that has been correctly identified; the second column is the total number of times that the gesture has been carried out. As it can be seen in Fig. 6, the hand recognition gesture works well for a 98% of the cases.

6 Conclusions

In this paper we have presented a real-time algorithm to track and recognise hand gestures for human-computer interaction within the context of videogames. We have proposed an algorithm based on skin colour hand segmentation and tracking for gesture recognition from extracted hand morphological features. The system's performance evaluation results have shown that the users can substitute traditional interaction metaphors with this low-cost interface.

The experiments have confirmed that continuous training of the users results in higher skills and, thus, better performances. Also the system has been tested in indoor laboratory with changing background scenario and low light conditions. In these cases the systems run well, with the logical exception of similar skin background situations or several hands intersecting in the same space and time. The system must be improved to discard bad classifications situations due to the segmentation procedure. But, in this case, the user can restart the system only going to the Start hand state.

Acknowledgements

The projects TIC2003-0931 and TIC2002-10743-E of MCYT Spanish Government and the European Project HUMODAN 2001-32202 from UE V Program-IST have subsidized this work. J.Varona acknowledges the support of a Ramon y Cajal fellowship from the spanish MEC.

References

- [1] V.I. Pavlovic, R. Sharma, T.S Huang, "Visual interpretation of hand gestures for human-computer interaction: a review", *IEEE Pattern Analysis and Machine Intelligence*, 19(7): 677 – 695, 1997 .
- [2] R. Bowden, D. Windridge, T. Kadir, A. Zisserman, M. Brady, "A Linguistic Feature Vector for the Visual Interpretation of Sign Language", in Tomas Pajdla, Jiri Matas (Eds), *Proc. European Conference on Computer Vision, ECCV04*, v. 1: 391-401, LNCS3022, Springer-Verlag, 2004.
- [3] J. Segen, S. Kumar, "Shadow gestures: 3D hand pose estimation using a single camera", *Proc. of the Computer Vision and Pattern Recognition Conference, CVPR99*, v. 1: 485, 1999.
- [4] M. Isard, A. Blake, "ICONDENSATION: Unifying low-level and high-level tracking in a stochastic framework", *Proc. European Conference on Computer Vision, ECCV98*, pp. 893-908, 1998.
- [5] C. Shan, Y. Wei, T. Tan, F.Ojardias , "Real time hand tracking by combining particle filtering and mean shift", *Proc. Sixth IEEE Automatic Face and Gesture Recognition, FG04*, pp: 229-674, 2004.
- [6] L. Bretzner, I. Laptev, T. Lindeberg , "Hand Gesture Recognition using Multi-Scale Colour Features, Hierarchical Models and Particle filtering", *Proc. Fifth IEEE International Conference on Automatic Face and Gesture Recognition, FRG02*, 2002 IEEE.
- [7] K.Ogawara, K. Hashimoto, J. Takamtsu, K. Ikeuchi, "Grasp Recognition using a 3D Articulated Model and Infrared Images", Institute of Industrial Science,. Univ. of Tokyo, Tokyo, Japan.
- [8] T. Heap, D. Hogg, "Wormholes in shape space: tracking through discontinuous changes in shape", *Proc. Sixth International Conference on Computer Vision, ICCV98*, pp. 344-349, 1998.
- [9] G.R. Bradski, "Computer video face tracking for use in a perceptual user interface," *Intel Technology Journal*, Q2'98, 1998.
- [10] D. Comaniciu, V. Ramesh, "Robust detection and tracking of human faces with an active camera" *Proc. of the Third IEEE International Workshop on Visual Surveillance*, pp: 11-18, 2000.
- [11] C.M. Bishop, *Neural Networks for Pattern Recognition*. Clarendon Press, 1995.
- [12] J. Varona, J.M. Buades, F.J. Perales, "Hands and face tracking for VR applications", *Computers & Graphics*, 29(2):179-187, 2005.
- [13] G.R. Bradski, V. Pisarevsky, "Intel's Computer Vision Library", *Proc of IEEE Conference on Computer Vision and Pattern Recognition, CVPR00*, v. 2: 796-797, 2000.